

Working Title: A Comparison of Approaches for Unplanned Sample Sizes in Phase II
Clinical Trials

By

Molly Olson

Thesis

Submitted to the Faculty of the
Graduate School of Vanderbilt University
in partial fulfillment of the requirements
for the degree of

MASTER OF SCIENCE

in

Biostatistics

May, 2017

Nashville, Tennessee

Approved (in progress):

Tatsuki Koyama , Ph.D.

Jeffrey Blume , Ph.D.

ACKNOWLEDGMENTS

TABLE OF CONTENTS

	Page
ACKNOWLEDGMENTS	ii
LIST OF TABLES	iv
LIST OF FIGURES	v
ABSTRACT	vi
Chapter	
1 Introduction	1
2 Background	3
3 Unplanned Sample Sizes	5
3.1 Unplanned (need word other than unplanned?) Sample Sizes in Second Stage	5
3.2 Unplanned (new word?) Sample Sizes in First Stage	6
3.2.1 <i>Chang et al.</i> Alternative Designs and Adaptation	6
3.2.2 Likelihood Design	8
4 Example	12
5 Results	13
5.1 Discussion	13
5.2 Appendix	13
5.3 Questions	14
5.4 Notes	14
REFERENCES	16

LIST OF TABLES

Table

Page

LIST OF FIGURES

Figure

Page

ABSTRACT

In this thesis, we develop

Chapter 1

Introduction

The introduction will talk about the motivation for the thesis. Introduce some examples of when we would need such methods.

Oncology phase II clinical trials are often used to evaluate the initial effect of a new regimen to determine if to warrant further study in a phase III clinical trial [1, 2, 3]. Simon's two-stage design [2] is a commonly used design in specifying sample sizes and critical values in phase II oncology clinical trials. Koyama and Chen [3] point out that it is common for actual sample sizes of these phase II trials to differ than the planned, pre-specified sample sizes. This could happen because of unanticipated accrument speed, drop-out rates are unexpected, and often multi-center trials can be slow in sharing information ... Currently, when attained sample sizes differ from planned, call these unplanned sample sizes, it is common practice to treat the attained sample sizes as planned. Though, when achieved sample sizes differ from planned, hypothesis testing using the attained sample sizes as if they were planned is not valid and hypothesis testing in these cases is not straightforward [1, 3]. Because of these reasons, extensions of Simon's design for hypothesis testing with unplanned sample sizes is important.

Talk about examples here

There have been many attempts to develop Frequentist methods that handle unplanned sample sizes in the second stage while using the planned stage I sample size, but my literature review found that there were only few Frequentist methods to handle unplanned sample sizes in both stage I and stage II. Likelihood based designs, that are able to be an extension of Simon's design, offer a nice solution to this problem because these designs offer flexi-

bility in sample size without inflation of type I error. In this paper, we discuss the different methods for Simon's design when the attained stage II sample size is different than planned and when attained sample sizes in both stages are different than planned. In chapter 4, we review a concrete example from a Likelihood-based clinical trial, and in chapter 5, results of a numerical and theoretical study comparing the Frequentist properties of approaches in the setting where both stages differ **wording** in different settings are presented.

Chapter 2

Background

Simon's design will go here. This section will also talk about extending/shortening a trial (unplanned sample sizes) and how recalculating as if it were the planned design will introduce bias and inflate type I error. Talk about prespecifying (maybe here?)

We will only talk about extensions to Simon's design, hypothesis testing, and only stopping for futility in this paper.

Simon's II stage designs for clinical trials are common designs for phase II oncology clinical trials [2]. In Simon's designs, the null hypothesis $H_0 : p \leq p_0$ is tested against the alternative $H_1 : p > p_0$, where p is the true response probability, p_0 is the highest probability of response that would indicate that the research regimen is uninteresting and p_1 is the lowest probability of response that would indicate that the research regimen warrants further investigation. Under these hypotheses, it is required that the type I error rate remain less than α and power remain above $1 - \beta$. The general framework of Simon's design includes a sample size and critical value in each of the two stages. Let n_1 denote the first stage sample size, n_t the sample size at the end of the second stage, r_1 the first stage critical value, and r_t the critical value for the end of the second stage. Let X_1 be the number of successes observed in the first stage and X_2 be the number of additional success in the second stage. In the first stage, n_1 patients are enrolled. If r_1 or fewer patients ($X_1 \leq r_1$) are successes, then the regimen is rejected and the trial is stopped for futility. If $r_1 + 1$ patients are successful, then the trial continues to the second stage. In the second stage, $n_2 = n_t - n_1$ patients are enrolled. If r_t or fewer out of the n_t patients are successful ($X_t = X_1 + X_2 \leq r_t$), the treatment is considered to be futile, otherwise if $r_t + 1$ patients succeed, the treatment is considered to be effective and warrant further study.

Design characteristics: Let B denote the cumulative binomial distribution function and b denote the binomial probability mass function. The probability of early termination with probability p in Simon's designs is given by $PET = B(r_1, p, n_1)$. The expected sample size for probability p is then $EN = n_1 + (1 - PET)n_2$. The probability of rejecting a drug for probability p is then $PR(p) = B(r_1, p, n_1) + \sum_{x=r_1+1}^{\min[n_1, r]} b(x, p, n_1)B(r-x, p, n_2)$. It is required that $PR(p) \geq 1 - \alpha$ and $PR(p) \leq \beta$. Given these constraints, it follows that unconditional conditional power, $UCP(p)$, given probability p , is given by $1 - PR(p) = 1 - \left(B(r_1, p, n_1) + \sum_{x=r_1+1}^{\min[n_1, r]} b(x, p, n_1)B(r-x, p, n_2) \right) = \sum_{r_1+1}^{n_1} \left\{ \sum_{x_2=r_t-x_1+1}^{n_2} b(x_2, p, n_2) \right\} b(x_1, p, n_1)$, and $UCP(p_1) \geq 1 - \beta$ and $UCP(p_0) \leq \alpha$.

Simon introduced optimal and minimax designs. An optimal two-stage design is a Simon's design in which minimizes the expected sample size under the null hypothesis, response value p_0 , (EN_0) while still satisfying the type I and type II error probability restrictions. The minimax design will minimize the maximum sample size ($n_1 + n_2$). Jung *et al.* [4] introduced an extension of Simon's designs called admissible designs that are considered a compromise between optimal and minimax designs because they have similar maximum sample size as the minimax design and a similar EN_0 to the optimal design. Admissible designs optimize a straight line on the (n, EN) -plane, $q \times n + (1 - q) \times EN$, for some $q \in [0, 1]$ [4]. Admissible designs satisfy (α, β) constraints and obtain an expected sample size somewhere between optimal and minimax designs. Admissible designs may be attractive because they have agreeable properties of both the minimax and optimal design. Simon does not allow for early termination of the trial for efficacy [2], and we do not consider that design here.

- We consider only redesigns here -something you can prespecify. Not calculate after you get the values.
- Talk about this in the "why we care"

Chapter 3

Unplanned Sample Sizes

This chapter will talk about unplanned sample sizes when only the second stage is different and when both stages can be different. The former will talk about methods such as Koyama and Chen, UMVUE, MLE, etc. The latter will talk about the likelihood design, Chang, adaptation of Chang, and possibly Wu.

3.1 Unplanned (need word other than unplanned?) Sample Sizes in Second Stage

When over- or under-enrollment occurs, a straightforward solution is to perform an interim analysis on the planned number of first stage subjects, and adjust the testing procedure for a sample size in the second stage that is different than planned. This is possible under the assumption of non-informative dropouts; stage one is concluded when the number of non-missing patients is equal to the planned stage one sample size, and if over enrollment occurs in the first stage, they will only be considered for the second stage analysis [3]. Literature exists describing point estimation of the response rate and p-values for hypothesis testing when stage two sample size is modified. A review of these methods can be found by Porcher et al. [1], though we only consider hypothesis testing when the trial design can be prespecified, rather than calculating a p-value from the attained data **call this proper p-value?**, here (I basically mean that we want to be able to prespecify the design with critical values for potential sample size deviations. A p-value is calculated after the trial and basically tries to get the correct p-value (can't use a conventional p-value) for trial deviations rather than pre-specifying rejection criteria/prespecifying a design

with those sample sizes). Koyama et al. propose a method for inference when stage 2 sample sizes deviate from the planned stage 2 sample size [3]. Let $n_1, n_t, r_1, r_t, \alpha$ and β be the original design parameters. They calculate a new critical value, r_t^* , by finding the value of r_t^* such that $P_{p_0}[X'_2 \geq r_t^* | X_1 = x_1] \leq P_{p_0}[X_2 \geq r_t(x_1) | X_1 = x_1] \dots$

3.2 Unplanned (new word?) Sample Sizes in First Stage

Decide later if subsections are needed.

Because accrument of patients can often be unexpected in the first stage, it's imperative that methods are available to handle situations with attained sample sizes that differ from the planned sample size. Green and Dahlberg [5] and Chen and Ng [6] propose methods for inference when first stage sample sizes differ than planned. Green and Dahlberg extended Southwest Oncology Group's inference method by suggesting to perform a hypothesis test on $H_0 : p = p_1$ versus $H_1 : p < p_1$ in the first stage at the 0.02 α -level and concluding futility if the p-value is ≤ 0.02 . They then suggest testing $H_0 : p = p_0$ versus $H_1 : p < p_0$ in the second stage at the 0.05 level. Li et al. indicate that this approach controls type I error and acheives desired power, though this approach is founded on an overall α -level of 0.05, and it is unclear how this method would generalize to any α -level [7]. Chang et al. also point out that Green and Dahlberg's designs can possibly be quite different than the planned designed. Chen and Ng —does this stuff— [6]. This is why this approach is not desired (I think you can prespecify this design though.)

3.2.1 *Chang et al.* Alternative Designs and Adaptation

Chang *et al* [8] proposed an alternative design that is an extension of Simon's two stage design in order to handle unplanned sample sizes in both the first and second stages. This method calculates new critical values for attained sample sizes, and thus one is able to

create and pre-specify a new design based on a preferred Simon or Admissible design in defense of the events of unplanned sample sizes (**basically trying to say in order to be ready with adjusted designs for unplanned sample sizes in case they occur**). Because it's desired to stay as closely to the original design as possible, we investigate this method using only attained first stage sample sizes while maintaining the original second stage sample size or original total sample size. Again, let n_1 , n_t , r_1 , r_t , p_0 , p_1 , α , and β be the original, planned design parameters. In the case that we let the total sample size be planned, let n_1^{**} be the attained sample size in the first stage and $n_2^{**} = n_t - n_1^{**}$. In the case that we let the second stage sample size remain as planned, let n_1^{**} again be the attained sample size in the first stage and $n_t^{**} = n_2 + n_1^{**}$.

Chang *et al* proposes that type II error probability spent in stage I, based on planned and attained sample size, is given by $\beta_1 = P(X_1 \leq r_1 | n_1, p = p_1)$ Based on the attained sample sizes, we choose to spend type II error in the first stage based on the the type II error probability spending function

$$\beta(m) = \begin{cases} \beta_1 m / n_1 & \text{if } m \leq n_1 \\ \beta_1 + (\beta - \beta_1)(m - n_1) / n_2 & \text{if } m > n_1 \end{cases}$$

We then find a new stage one critical value, s_1 , based on this probability spending function such that $P(X_1 \leq s_1 | n_1^{**}) \approx \beta(n_1^{**})$, where \approx means “closest to.” After s_1 is selected, we then search for an integer for the second stage critical value, s_t , that satisfies

$$\begin{aligned} & P(X_1 > s_1, X_t > s_t | n_1^{**}, m_2, p_0) \\ &= \sum_{s_1}^{n_1^{**}} P(X_2 > s_t - X_1 | X_1 = x_1) P(X_1 > s_1) \\ &\leq \alpha \end{aligned}$$

where $m_2 = n_2$ or n_2^{**} . Chang et al.'s design can be used for any α -level and are flexible, close to the original design, and preserve desired Frequentist characteristics.

Because we prefer to be conservative when straying from a desired Simon or Admissible design, we modify Chang et al.'s design by selecting s_1 that preserves the probability of early termination under the null. We select s_1 such that

$$P(X_1 \leq s_1 | n_1^{**}, p_0) \approx P(X_1 \leq r_1 | n_1, p_0)$$

We then select the stage two critical value, s_t , in the same fashion as Chang's design.

- Read chen and ng....

3.2.2 Likelihood Design

Briefly, the likelihood stage II design uses the likelihood ratio, as opposed to a p-value, as a measure of evidence [9]. Here, the likelihood ratio is

$$\begin{aligned} \text{LR}_n &= \frac{L_n(p_1)}{L_n(p_0)} \\ &= \frac{p_1^{x_t} (1 - p_1)^{n_t - x_t}}{p_0^{x_t} (1 - p_0)^{n_t - x_t}} \\ &\in \{[0, 1/k], [1/k, k], [k, \infty)\} \end{aligned}$$

and has three evidential zones: evidence for the null hypothesis, weak evidence, and evidence for the alternative hypothesis. If the $\text{LR}_n \in [0, 1/k]$, there is evidence for the null hypothesis, if $\text{LR}_n \in [1/k, k]$, there is weak evidence for either hypothesis, and if $\text{LR}_n \in [k, \infty]$, there is evidence for the alternative hypothesis.

Likelihood characteristics:

$$LR_n = \frac{L_n(\theta_1)}{L_n(\theta_0)} \in \{[0, 1/k], [1/k, k], [k, \infty)\}$$

Probability of Weak Evidence

$$\gamma_p = P(k_a \leq LR_n \leq K_b | H_p), k_a \leq k_b$$

Probability of Strong Evidence

$$\eta_1 = P(LR_n > k_b | H_1)$$

$$\eta_1 = P(LR_n < k_a | H_0)$$

Probability of Observing Misleading Evidence

$$\tau_0 = P(LR_n > k_b | H_0), \tau_0 \leq 1/k_b$$

$$\tau_1 = P(LR_n < k_a | H_1), \tau_1 \leq k_a$$

$$\tau_i = O_p(n^{-1/2}) \text{ instead of remaining fixed like type I error}$$

Translating likelihood properties into Simon-like design:

Interim: Translating to successes. This is the region in which we move to stage 2

$$UB_{interim} = \frac{\log(k_{bi}) - n_1 \log(\frac{1-p_1}{1-p_0})}{\log(\frac{p_1(1-p_0)}{p_0(1-p_1)})}$$
$$LB_{interim} = \frac{\log(k_{ai}) - n_1 \log(\frac{1-p_1}{1-p_0})}{\log(\frac{p_1(1-p_0)}{p_0(1-p_1)})}$$

(LB, UB) is the interval for weak evidence. If this was Simon's design, $LB_{interim} = r_1$

Probability of strong, misleading, and weak evidence under the null

$$P(\text{Strong}_{0i}) = B(\lfloor LB_{interim} \rfloor, n_1, p_0)$$

$$P(\text{Misleading}_{0i}) = 1 - B(\lfloor UB_{interim} \rfloor, n_1, p_0)$$

$$P(\text{Weak}_{0i}) = B(\lfloor UB_{interim} \rfloor, n_1, p_0) - B(\lfloor LB_{interim} \rfloor, n_1, p_0)$$

Probability of strong, misleading, and weak evidence under the alternative

$$P(\text{Strong}_{1i}) = 1 - B(\lfloor UB_{interim} \rfloor, n_1, p_1)$$

$$P(\text{Misleading}_{1i}) = B(\lfloor LB_{interim} \rfloor, n_1, p_1)$$

$$P(\text{Weak}_{1i}) = B(\lfloor UB_{interim} \rfloor, n_1, p_1) - B(\lfloor LB_{interim} \rfloor, n_1, p_1)$$

note: under Simon's, PET = 1-P(Weak)

Translating likelihood properties into Simon-like design:

Final Stage: Translating to successes.

The amount of successes that allow for continuation to the second stage are:

$$(\lfloor LB_{interim} + 1 \rfloor, \lfloor \min(n_1, UB_{interim}) \rfloor)$$

Probability of strong, misleading, and weak evidence under H_p

$$\begin{aligned} P(Weak_p) &= \sum_{x=\lfloor LB_{interim}+1 \rfloor}^{\lfloor \min(n_1, UB_{interim}) \rfloor} \left(b(x, n_1, p_p) \times B(UB_{interim} - x, n - n_1, p_p) \right) - B(LB_{interim} - x, n - n_1, p_p) \\ P(Strong_p) &= P(Strong_{0i}) + \sum_{x=\lfloor LB_{interim}+1 \rfloor}^{\lfloor \min(n_1, UB_{interim}) \rfloor} \left(b(x, n_1, p_0) \times B(LB_{interim} - x, n - n_1, p_0) \right) \\ P(Misleading_p) &= P(Misleading_{0i}) + \sum_{x=\lfloor LB_{interim}+1 \rfloor}^{\lfloor \min(n_1, UB_{interim}) \rfloor} \left(b(x, n_1, p_p) \times (1 - B(UB_{interim} - x, n - n_1, p_p)) \right) \end{aligned}$$

If we want to translate likelihood design into a Simon's design, we overwrite the LR limits

above as:

$$\begin{aligned} k_{ai} &= OR^{r_1} \frac{1 - p_1^{n_1}}{1 - p_0} = \frac{1 - p_0^{r_1 - n_1}}{1 - p_1} \frac{p_1^{r_1}}{p_0} \\ k_a &= OR^r \frac{1 - p_1^n}{1 - p_0} = \frac{1 - p_0^{r - n}}{1 - p_1} \frac{p_1^r}{p_0} \\ k_{bi} &= k_b = \infty \end{aligned}$$

Chapter 4

Example

Here put the results of comparing the Chang et al paper and adaptation to the protocol of the study. We used Monte Carlo simulation to examine the performance of the study design of Chang et al...

Chapter 5

Results

The findings will be put here - primarily tables covering different combinations of Simon's designs and unplanned sample sizes.

Should maybe talk about how I couldn't replicate some results in the paper.

Some differences with the likelihood:

controlling type 1 error, but criteria is controlling PET - I think T1E will still be controlled.

Assuming stage II sample size and R^2 is the same. We can add cohorts at the end of stage II. Talk about this.

5.1 Discussion

- Numerical study shows that keeping n the same as the original has better properties.
- Talk about when the adaptation design and chang design will differ. (extreme sample size shifts? I cant remember.)
- Compare the design approaches

5.2 Appendix

Maybe put the chang design here that goes beyond stage 2 sample size being origianl total and original second stage?

5.3 Questions

- Particularly when describing other peoples' methods, how do you cite?
- If likelihood design can translate into simon's design and have better properties - can that be used as a tool to come up with these designs instead of using Chang's method? Or are they different because you're operating under frequentist vs likelihood inference?
- Chang's method - would we use a conventional p-value?
- Koyama's method: equation (5) does X_2' have to be attained? is this calculated after stage 2 happened? If so, isn't that not what we are focusing on? aren't we trying to prespecify a design rather than deal with what happened?
- Why do we choose PET closest to the original? - consider choosing closest PET, being conservative. mention that there are two ways to think about it. one is conservative and one is closest and they're pretty close to each other.

5.4 Notes

- So I think a lot of what's going on is they are calculating p-values based on the attained sample size, but not actually a second stage critical value.
- Are papers like Tatsuki's finding critical values? Or just calculating p-values? Tatsuki calculates a new critical value.
- conclusions: read green and dahlberg, chen and ng, tatsukis, and zhao. I'm fairly certain tatsuki's is the only one with a new critical value when stage 2 is different. GD, CN may be "redesigns" The bayesian paper has a good explanation of faults in green and dahlberg, chen and ng. we can talk about these briefly in the first stage differing.

- inference from likelihood is more straightforward. the authors in koyama have shown that the p-value depends on the design. p-value's in multistage designs depend on the p-value. first stage is more complex.

REFERENCES

- [1] R. Porcher and K. Desseaux, “What inference for two stage phase II trials?,” *BMC Medical Research Methodology*, vol. 12, p. 117, 2012.
- [2] R. Simon, “Optimal two-stage designs for phase ii clinical trials,” *Controlled Clinical Trials*, vol. 10, pp. 1–10, 1989.
- [3] T. Koyama and H. Chen, “Proper inference from simon’s two-stage designs,” *Statistics In Medicine*, vol. 27, pp. 3145–3154, 2008.
- [4] S.-H. Jung, Y. Lee, K. Kim, and S. L. George, “Admissible two-stage designs for phase ii cancer clinical trials,” *Statistics In Medicine*, vol. 23, pp. 561–569, 2004.
- [5] “Planned versus attained design in phase ii clinical trials,” *Statistics in Medicine*, vol. 11, pp. 853–862, 1992.
- [6] T. Chen and T. Ng, “Optimal flexible designs in phase ii clinical trials,” *Statistics in Medicine*, vol. 17, pp. 2301–2312, 1998.
- [7] “A bayesian approach for unplanned sample sizes in phase ii clinical trials,” *Clinical Trials*, vol. 9, pp. 293–302, 2012.
- [8]
- [9] J. Blume and D. Ayers